

SYSTEM AND METHOD FOR HANDLING FLOWS IN A NETWORK

Claim to Priority

This application claims the benefit of United States Provisional Application No. 60/149,174, filed on August 17, 1999.

Field of the Invention

The field of the invention is handling flows in a network, and in particular handling packets that relate to the same conversation as a part of a flow.

Brief Description of the Drawings

Figure 1 shows a switch that handles a flow between two hosts in accordance with an embodiment of the present invention.

Figure 2 shows a flow that passes through one Ethernet switch between two hosts in accordance with an embodiment of the present invention.

Figure 3 shows flows between two switches and two hosts in accordance with an embodiment of the present invention.

Figure 4 shows multicast flows in accordance with an embodiment of the present invention.

Summary of the Invention

A flow in a network is identified and handled by using a virtual host address. A packet is received at a switch with a first virtual host address as its destination address. If the packet is the first packet of a flow received by the switch, then a second virtual host address is determined by the switch. The first virtual host address is stored in a packet forwarding table correlated with the second virtual host address. A subsequently received packet of the same flow has the same first virtual host address as its destination address, and is forwarded to the second virtual host address in accordance with the packet forwarding table.

Detailed Description

The wide area network is evolving to one that integrates virtual circuit switching (label

swapping) for flows with conventional datagram forwarding. A first step along that road was described by Ipsilon by Newman, P et al, in IP Switching - ATM Under IP, IEEE Trans. on Networking, Vol .6, No.2, April 1998, which:

- a) uses a classification algorithm to detect flows among the influx of IP packets;
- b) uses IP datagram forwarding to determine where to send the packet;
- c) creates a virtual circuit connection through the switch to the same place that the IP packet is being sent;
- d) transmits the VCI of that connection to the upstream switch with an indication that subsequent packets should be encapsulated with that VCI; and
- e) arranges that incoming packets encapsulated with that VCI are switched not routed.

We have modified this concept to provide flow switching on local area networks (LANs) that use Ethernet. Figure 1 illustrates a switch that handles a flow between two hosts, H and K. Usually, Ethernet addresses are of hosts rather than endpoints of flows. Our design uses Ethernet addresses to also identify flows on the LAN. It is exactly as if the switch contains one virtual host for every flow. The Ethernet address of that virtual host, referred to here as V, is temporarily assigned from a block of locally administered Ethernet addresses. Packets of a flow from host H to host K pass through the virtual host V. The source and destination addresses in packets leaving H are H and V respectively. Packets traveling from V to K have source and destination addresses equal to V and K. The switch performs Ethernet address swapping as follows:

- a) the destination address of an incoming packet is moved into the source address field; and
- b) a new destination address is obtained from a "VC forwarding table" held within the switch.

The technique is compatible with existing applications of Ethernet because in effect all we have done is to add extra (virtual) hosts to the network.

Whereas the Epsilon technique used a classification algorithm to detect flows among IP packets, we have experimented with the idea that the host application should make that decision. We have added a single byte, `vc_flag`, in the general socket structure of our hosts to say that the application wants special service for the flow of packets passing through the socket. The presence of that flag tells the socket software to use a virtual host Ethernet address instead of the destination Ethernet address implied by the IP header.

The switch does traditional Ethernet packet forwarding on all packets except those that are addressed to a virtual host. Packets addressed to a virtual host are switched using data in a VC forwarding table. The first packet for a new flow causes an entry to be made in the VC forwarding table based upon the IP destination contained in the packet.

By this means we have created in the local area a sufficient means to provide quality communication service on a per-flow basis. When the technique is matched to flow switching in a wide area network the user has full benefit of end-to-end flow switching, from a socket in one host to a socket in another. This has been achieved with minimal impact on host software, no interference with existing applications, and complete compatibility with existing Ethernets.

Ethernet RFC 894 packet format

DESTINATION ADDRESS
SOURCE ADDRESS
TYPE
PAYLOAD
PAD
FRAME CHECK

The packet forwarding table used by each virtual host is constructed by examining the header of the first IP packet in a flow.

Of course, virtual hosts do not really exist, even as processes within a switch. It is just that the actions of a switch as seen from outside are exactly as described by the model. Internally the switch uses a combination of technologies found today in IP routers and virtual circuit switches. It is a table-driven process that stores packets in queues, processes their headers and transfers them to the appropriate output ports with appropriate attention to the quality of service appropriate to each traffic class.

The same technique can be used for point to multipoint flows, as shown in Figure 4. In this example, host H is the root of a multicast tree that transmits packets to the two hosts K and L. The forwarding table now has three rows, one for each host in the multicast, and a third column indicates which host is the "root" of the multicast tree. Packets coming from H are copied to each of the hosts given in the other rows of the table. Packets addressed to V from K and L may either be rejected or propagated upstream depending upon the permission stated in the "perm" column. Note that if K and L do transmit packets upstream, H must examine the IP header to determine the source of each packet.

An example of a virtual circuit signaling connection set-up protocol follows. A protocol for setting up a connection between two hosts A and B takes place in three stages. First A requests that the connection be made, then B accepts the request and causes a virtual circuit to be created, and finally A confirms that indeed there is a connection.

The connection request is sent as an ordinary IP datagram from A to B. The accept message is sent as a signal, which is a message from A to B that is flagged for special attention in each of the network nodes along the way. As this signal progresses through the network a (full duplex) virtual circuit is created between A and B. Finally, the confirmation message from A is transmitted over the new virtual circuit.

